

JOHNSON GRANT
IN-63-CR

7576

(NASA-CR-188091) QUEUEING MODELS FOR TOKEN
AND SLOTTED RING NETWORKS Thesis (Houston
Univ.) 35 p CSCL 09B

N91-21791

Unclas
G3/63 0007596

QUEUEING MODELS FOR TOKEN AND SLOTTED RING NETWORKS

Dissertation Proposal

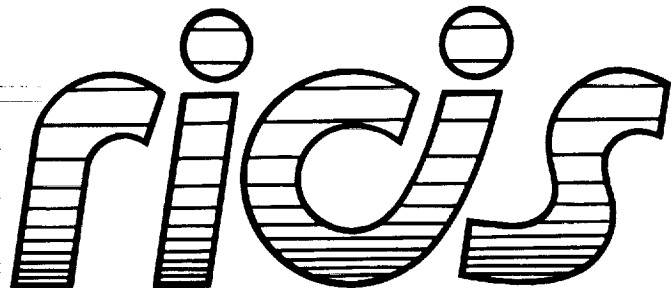
Jeffery H. Peden

Digital Technology

September 1990

**Cooperative Agreement NCC 9-16
Research Activity No. SE.31**

**NASA Johnson Space Center
Engineering Directorate
Flight Data Systems Division**



**Research Institute for Computing and Information Systems
University of Houston - Clear Lake**

T · E · C · H · N · I · C · A · L R · E · P · O · R · T

The RICIS Concept

The University of Houston-Clear Lake established the Research Institute for Computing and Information systems in 1986 to encourage NASA Johnson Space Center and local industry to actively support research in the computing and information sciences. As part of this endeavor, UH-Clear Lake proposed a partnership with JSC to jointly define and manage an integrated program of research in advanced data processing technology needed for JSC's main missions, including administrative, engineering and science responsibilities. JSC agreed and entered into a three-year cooperative agreement with UH-Clear Lake beginning in May, 1986, to jointly plan and execute such research through RICIS. Additionally, under Cooperative Agreement NCC 9-16, computing and educational facilities are shared by the two institutions to conduct the research.

The mission of RICIS is to conduct, coordinate and disseminate research on computing and information systems among researchers, sponsors and users from UH-Clear Lake, NASA/JSC, and other research organizations. Within UH-Clear Lake, the mission is being implemented through interdisciplinary involvement of faculty and students from each of the four schools: Business, Education, Human Sciences and Humanities, and Natural and Applied Sciences.

Other research organizations are involved via the "gateway" concept. UH-Clear Lake establishes relationships with other universities and research organizations, having common research interests, to provide additional sources of expertise to conduct needed research.

A major role of RICIS is to find the best match of sponsors, researchers and research objectives to advance knowledge in the computing and information sciences. Working jointly with NASA/JSC, RICIS advises on research needs, recommends principals for conducting the research, provides technical and administrative support to coordinate the research, and integrates technical results into the cooperative goals of UH-Clear Lake and NASA/JSC.

***QUEUEING MODELS FOR TOKEN
AND
SLOTTED RING NETWORKS***
Dissertation Proposal

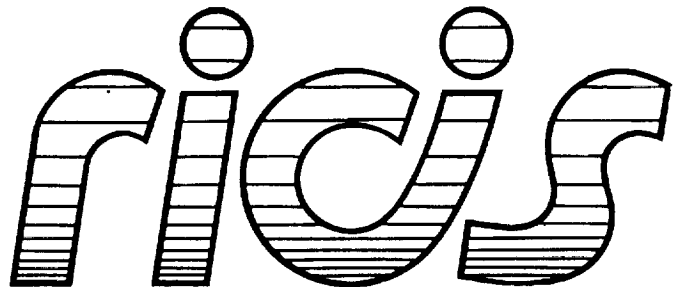
Jeffery H. Peden

Digital Technology

September 1990

Cooperative Agreement NCC 9-16
Research Activity No. SE.31

NASA Johnson Space Center
Engineering Directorate
Flight Data Systems Division



*Research Institute for Computing and Information Systems
University of Houston - Clear Lake*

T · E · C · H · N · I · C · A · L R · E · P · O · R · T

Preface

This research was conducted under auspices of the Research Institute for Computing and Information Systems by Jeffery H. Peden and Digital Technology. Dr. George Collins, Associate Professor of Computer Systems Design, served as RICIS technical representative for this activity.

Funding has been provided by the Engineering Directorate, NASA/JSC through Cooperative Agreement NCC 9-16 between NASA Johnson Space Center and the University of Houston-Clear Lake. The NASA technical monitor for this activity was Frank W. Miller, of the Systems Development Branch, Flight Data Systems Division, Engineering Directorate, NASA/JSC.

The views and conclusions contained in this report are those of the author and should not be interpreted as representative of the official policies, either express or implied, of NASA or the United States Government.

Queueing Models for Token and Slotted Ring Networks

Dissertation Proposal

Jeffery H. Peden

1 Introduction

Currently the end-to-end delay characteristics of very high speed local area networks are not well understood. As networks begin to be used for real-time and near real-time applications, the entire delay incurred by a packet becomes important, due to the fact that packet processing time can completely subsume transmission time. It is also important that queue length characteristics be understood so that adequate buffer space can be allocated in the network.

The transmission speed of computer networks is increasing, and local area networks especially are finding increasing use in real time systems. Therefore, in order to model accurately total network delay, queueing models must include all characteristics of packet processing. Since often only a small fraction of total delay is due to actual transmission, models which deal only with the MAC (Medium Access Control) layer are becoming increasingly inadequate.

Ring network operation is generally well understood for both token rings and slotted rings. Many models exist for several MAC layer service disciplines, for example, exhaustive service, gated service, and single packet per token service. There is, however, a severe lack of queueing models for higher layer operation.

There are several factors which contribute to the processing delay of a packet, as opposed to the transmission delay. These are the packet's priority, its length, the user load, the processor load, the use of priority preemption, the use of preemption at packet reception, the number of processors, the number of protocol processing layers, the speed of each processor, and queue length limitations. Any useful queueing model must take all these factors into account.

1.1 The ISO OSI Protocol Stack

The currently accepted reference model for packet processing prior to and after transmission is the ISO OSI protocol stack. This is a set of conceptual functions which are helpful in ensuring the delivery of a packet across any network. The model is network independent since it only defines the necessary actions, not how they are implemented.

The OSI model consists of seven layers: the Application Layer, Presentation Layer, Session Layer, Transport Layer, Network Layer, Data Link Layer, and Physical Layer (see Figure 1). The Data Link Layer is further subdivided into the Logical Link Control and Medium Access Control sublayers. Each layer represents one or more conceptual functions that are performed by the network each time it transmits or receives a packet. These layers are conceptual in that they need not be implemented as distinct entities.

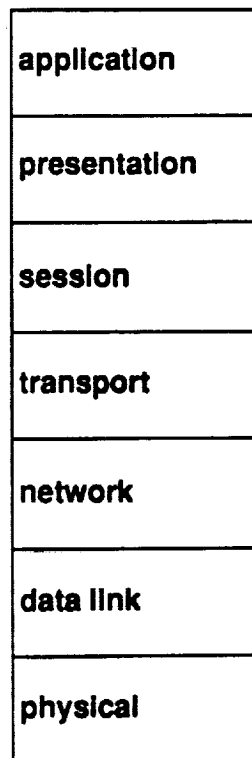


Figure 1

The application layer is the interface between the user's application and the functions provided by the network. This layer supplies such capabilities as file transfer and remote login. The application layer is the end user's view of the network, and is independent of the actual physical configuration, that is, for example, whether the network is an Ethernet, a token ring, or a wide area network. Each of the different services provided have different delay and queue length characteristics. For example, a file transfer would tend to submit fairly long packets to the network in a deterministic manner, whereas a remote login would probably cause much shorter packets to be submitted in a more random manner.

The presentation layer negotiates and controls transfer syntax. This is necessary since not all devices use identical syntax for the representation of data values. At present the presentation layer is especially poorly understood, and thus almost no attempts have been made to develop queueing models of this layer's functions.

The session layer sets up a dialogue between communicating devices. It ensures that any transactions which take place during communication are completed in their entirety, thus making communicating entities appear to be either completely operational or nonexistent.

The transport layer ensures the end-to-end delivery of all packets transmitted. This enables reliable transmission over the network, making the network appear as if it is a perfect channel. Another responsibility of the transport layer is the segmentation of messages into packets. This is often necessary since messages may be larger than the maximum size packet that the Network or Data Link layers can handle. The transport layer, and all layers above it, are end-to-end layers in that they have the illusion of communicating directly with their peer layers at the "receiving" end.

The network layer is responsible for packet routing. It also performs segmentation if differing network segments have different maximum packet sizes. The network layer is the highest layer in the OSI stack that is not end-to-end. That is, network layers at each node only talk to

routers on the same network segment; thus the network layer at the initiation end of a transmission does not talk to the network layer at the final reception end of a transmission.

The data link layer is responsible for node-to-node transmission of packets. This may or may not be done reliably, depending upon the network and the wishes of the user. The data link layer is not an end-to-end layer. This node-to-node transmission is accomplished through the use of the Logical Link Control sublayer, which initiates the transmission to and handles the reception from adjacent nodes, and the Medium Access Control sublayer, which mediates access to the actual transmission medium. The Logical Link Control is protocol independent, whereas the Medium Access Control is protocol dependent.

The physical layer is the mechanism that actually transmits bits onto the network medium. It performs data encoding/decoding, carrier sense, fault detection, etc. The physical layer attempts to isolate the higher layers of the protocol stack from the actual physical processes involved in transmission. Ideally, the higher layers in the protocol stack do not know, for example, whether the network is a token ring or an Ethernet.

1.2 Modeling the Protocol Stack

1.2.1 Inherent Problems

Protocol stack models can be developed in two basic forms: mathematical and simulation models. Since the ISO OSI model is intended as a reference model only, it is not a simple matter to write simulation code or to develop a parameterized mathematical model. Implementations of the OSI protocol stack almost always depart from a strict process or function-per-layer implementation method by combining or coalescing layers for the sake of efficiency in terms of memory, processing speed, the number of processors, or some combination of the above.

Other factors which must be taken into account when modeling the OSI protocol stack are the actual functions provided (e.g., the stack being modeled may not implement the full func-

tionality of the OSI reference model), the speed of the processors, the amount of memory available (available memory puts an upper bound on maximum queue size), maximum packet length (since certain layers perform segmentation), the change from packet arrivals to bulk arrivals at certain layers (if segmentation is used), whether or not separate processors are used for transmission and reception, and if not, which takes precedence, and expected error rates.

It is also necessary to account for the fact that certain layers in the stack may not be used for every packet transmission, or indeed for every message transmission. For example, the session layer is responsible for setting up the dialogue between two stations, but once this function has been performed, the session layer becomes largely inactive in its effects on transmission delay (except for maintaining checkpoints). That is, there is normally some initial one-time overhead involved in setting up a communications path between two entities. Therefore, the frequency of occurrence of this overhead will have an impact on overall delay.

1.2.2 Methods

Since the actual implementations of protocol stacks will vary both in architecture and functionality, any attempt at modeling must account for this by being parametric and as general as possible. Therefore, stack models derived here will not have a strict one-to-one mapping to the OSI reference model.

The technique used here will be to model stack processing as a collection of one or more queueing systems connected in series. We assume that the arrival process at the initial (and possibly only) queue is Poisson. If there are more queueing systems in the series, then each has what we term a k^{th} level *series* arrival process (abbreviated S_k). It is also possible that some of the queues in the series will have a k^{th} level *bulk series* arrival process (abbreviated S_k^B) due to segmentation in the previous queueing system. Figure 2 shows the general queueing system to be modeled.

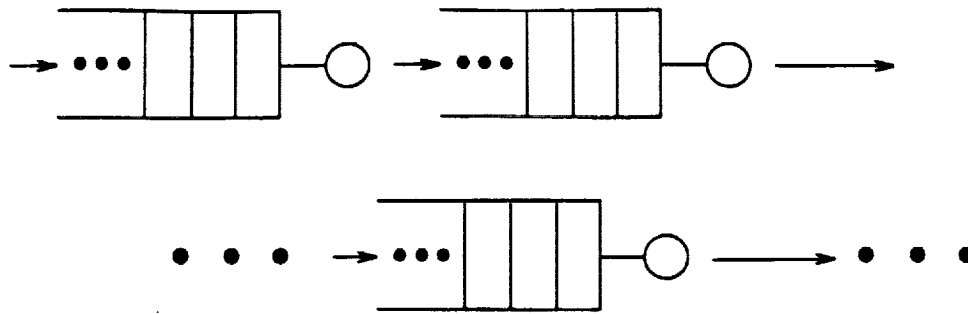


Figure 2

Once solutions have been found, it is then possible to map various configurations of the above system to actual implementations of protocol processing stacks. Note that it is not necessary that these implementations be of the ISO protocol stack; we can map the above system to other protocol stacks as well, such as a "short" ISO stack, or an implementation of the TCP/IP protocol.

2 Previous Work

This section summarizes existing work in the field relevant to this dissertation. Information covered deals with token and slotted rings, priority queues, upper layer modeling, and queues in series.

In their seminal paper on cyclic service [KOH72], Konheim and Meister present mathematical results dealing with exhaustive cyclic service. Their model is derived from first principles, and is an excellent introduction to the field. This paper is especially useful in that the mathematics are applicable (with suitable modifications) to other forms of cyclic service. The main drawbacks of this paper are that it does not deal with nonexhaustive cyclic service, and it neglects all mention of priority queueing.

Heyman [HEY83] takes many of the results from [KOH72] and applies them to an existing ring network. The usefulness of this paper lies in the fact that it is an excellent example of the application and modification of Konheim and Meister's results to an existing network. Its drawbacks are essentially those of [KOH72] in that it does not deal with exhaustive service or priority queueing.

In his paper on nonexhaustive cycle service [KUE79], Kuehn develops an analytic model for cyclic service using a single packet per token service discipline. In this paper, he derives formulas for the mean and variance of waiting time and queue lengths. He makes an important step in his derivation by taking into account the effects of each token cycle on the one following. The main drawback lies in his assumption that all queues except the one transmitting are in equilibrium, and are thus not affected by the transmission of the queue in question. This assumption has the effect of simplifying the mathematics, but it also causes the analysis to be at variance with actual network operation in many cases.

Levy and Yechiali produce an important result concerning queue lengths and delay in their paper [LEV75]. They prove that the delay of a packet in a cyclic service system can be decom-

posed into queueing time plus the time spent waiting for the token. This decomposition method has a simplifying effect on the analysis of cyclic service systems.

In their paper on token passing protocols [REGO84], Rego and Ni develop a method of dealing with the serial dependence of token cycles. They derive methods for computing the limit of dependence as well computing the covariance of the dependence.

Genter and Vastola develop an approximate model of gated limited service in [GENT88]. This service method allows a maximum number of packets to be served at token arrival, this number being the minimum of the number of packet present at token arrival and a preset maximum. This is an important result; however, the authors ignore the effects of service time on the token cycle time, as well as the effects of priority.

In their paper [APPL86], Appleton and Peterson derive the probability density function for a single packet per token service discipline. This paper actually uses a protocol definition as the basis for its analysis, which is a slight modification of the IEEE 802.5 token ring definition. This result is important as it is the only known paper to derive the probability density function for this type of service, the others being content to derive transform equations.

The correspondence paper by Everitt [EVER89] derives the pseudoconservation laws for both gated limited and exhaustive limited service without attempting to derive either transform equations or probability density functions for these service disciplines. This is nevertheless an important result as it gives certain stability criteria for ring functioning when these disciplines are used.

In two of my previous papers [PEDE87], [PEDE88], I derive approximate models for exhaustive limited priority service. These models are shown to be very accurate via comparison with simulation results. The drawbacks of these papers are that several simplifying assumptions are made. Furthermore, no transforms are derived; only mean value approximations are given. However, these results are important as a first step in the understanding of this form of service.

Another important feature is that certain relationships are shown governing the effects of a packet's priority with the load, number of ring stations, and packet size.

In their paper [TAKA86], Takagi and Murata make a queueing analysis of a non-preemptive reservation priority token ring network. Their analysis assumes single-packet-per-token service, which is the main drawback of the paper. The main strength of the paper is that it does not assume a specific number of priorities.

Manfield presents an analysis of two-way priority traffic for a polling system in his paper [MANF85]. His derivation assumes an interdependence between incoming and outgoing traffic, giving priority to outgoing traffic. This paper is important because of the above assumption, as this has an impact of upper layer performance in an actual network.

In her paper [GORU87], Gorur derives the equations governing the limiting throughput levels for the various ring priorities when a timed token ring access method is used. However, formulas for mean delay and queue length are not given. This is nevertheless an important result since equations governing operational bounds are derived.

In her dissertation [SCHU89], Schult derives a single packet per token analytic model for priority operation on a token ring. These results also take into account the effects of transmissions on successive token cycles. The main drawback to these results is that only single packet per token service is treated.

Zafirovic-Vukotic and Niemegeers present an excellent single priority treatment of a high speed slotted ring protocol [ZAFI87]. Mean value equations are derived for both delay and queue lengths, and the results of these equations are compared with simulation results. The main drawback of this paper is the lack of any priority queueing results.

In their paper on bandwidth allocation [LI88], Li and Zarki present an analysis on the dynamic allocation of bandwidth on a slotted ring. The importance of this paper is that it models different classes of users, thus having a direct application to priority service on a slotted ring.

network. The method of modeling used is traffic prediction.

Bhuyan, Ghosal and Yang present an approximate model for interconnected ring networks in their paper [BHUY89]. The ring networks studied are token rings, register insertion rings, and slotted rings. All the models are based on nonexhaustive service and potentially infinite queue lengths at each station.

The paper by Mills, Wheatley and Heatley [MILL87] is an introduction to the techniques and importance of modeling the upper layers in a protocol stack. They show that the delay induced by packet processing can completely subsume the delay experienced at the MAC layer. This lays the groundwork for other research in this area. The main drawback of this paper is that no analytic modeling is done; the entire paper is concerned with simulation systems.

In his master's thesis [STRA88], Strayer presents some very important results concerning the effects of segmentation on the delay and throughput of a packet. This thesis is concerned almost entirely with upper layer packet processing, and deals with actual implemented ring networks, not simulations. The main drawback of this source is the fact that most of the data was drawn from a single network, thus making it necessary to look at data trends rather than actual values.

Murata and Takagi develop a two-layer modeling technique for local area networks in [MURA88]. This paper describes a MAC layer submodel and a transport layer submodel. Both models deal with priority traffic and acknowledgements. An iterative technique is developed as a solution algorithm.

A different method of end-to-end modeling is proposed in the paper by Mitchell and Lide [MITC86]. Their model uses a hierarchical modeling technique and an iterative technique for solution which combines analytic and simulation methods.

My dissertation will unify the work described above, extend it to handle exhaustive limited (EL) service at the MAC layer, and multiple priority messages at all layers, and extend the model

to the upper ISO layers.

3 Intended Contributions

The intended contribution of this dissertation is to extend currently existing medium access queueing models by adding modeling techniques which will handle exhaustive limited service (see Section 4) both with and without priority traffic, and to extend modeling capabilities into the upper layers of the OSI model. There are certain limitations which must be taken into account, however, the main one being that all arrivals are assumed to be exponential in nature. If this assumption is not made, the problem becomes one of sequential G/G/m priority queues, for which only approximations are available.

Some of the models presented will be parameterized solution methods, since it will be shown in section 4 that certain models (for example, the single priority exhaustive limited service model) do not exist as parameterized solutions, but rather as solution methods. Since both token rings and slotted rings are dealt with, there will be two distinct solution methods presented for the medium access control queueing models.

3.1 Dissertation Outline

The proposed outline of the dissertation is as follows:

I. Introduction and explanation of network operation

- A. ring networks
 - 1. token rings
 - 2. slotted rings
- B. layers
 - 1. functions
 - 2. modeling

II. Review of existing network performance models

- A. token rings
 - 1. service types
 - 2. timed token
 - 3. reservation token
- B. slotted rings
 - 1. service types
 - a. release slot at reception
 - b. release slot at acknowledgement
- C. upper layers

1. priority preemption/no preemption
2. preemption/no preemption on packet reception

III. The exhaustive limited (EL) service model

- A. significance
- B. development
 1. initial M/G/1 model
 2. modification to S/G/1
- C. application to token ring
 1. timed token
 2. reservation token
- D. application to slotted rings

IV. The priority model

- A. significance
- B. development
- C. combining with the EL model
 1. how the model relates to the EL model
 2. development
- D. token rings
 1. timed token
 2. reservation token
- E. slotted rings

V. Modeling the upper layers

- A. development
 1. single priority
 2. multiple priority
- B. token rings
 1. timed token
 2. reservation token
- C. slotted rings

VI. The simulation

- A. capabilities
 1. networks modeled
 2. upper layer models
 3. classes
- B. limitations
 1. steady state load only
 2. no initial load, no trace mode
- C. statistical validation
 1. stability
 2. confidence intervals

VII. Comparing simulation to analytic models

- A. token rings
 1. timed token
 2. reservation token
 3. differing numbers of packets per token
 4. various priority levels

- B. slotted rings
 - 1. slot release at reception
 - 2. slot release at acknowledgement
- C. upper layers

VIII. Simulating class traffic and entry layer modification

- A. modifying classes
 - 1. equally divided loads
 - 2. skewed loads
- B. modifying entry layer
 - 1. all packets enter at the same layer
 - 2. different classes enter at different layers
- C. simulation of non-homogeneous loads and varying entry layers

IX. Application of the analytic and simulation models

- A. mapping the analytic model to actual implementations
 - 1. correspondence of layers to queues
 - 2. layer groupings
- B. mapping the simulation to actual implementations
 - 1. memory utilization
 - 2. practical design tradeoffs
- C. Using the analytic and/or simulation models for prediction data

X. Conclusions

- A. main conclusions
 - 1. EL model validity
 - 2. priority model validity
- B. future work
 - 1. mathematical modeling of class traffic
 - 2. modeling changing priorities and upper layers

3.2 Discussion

Chapter I is a brief introduction to the subject of ring networks. It covers the ring concept, token rings, and slotted rings, and explains briefly the abstractions which we will model. We also cover the basic idea of the protocol stack, and why it is important that we model this network feature.

Chapter II is a review of existing simulation and analytic models. It covers the relevant work concerning token rings, slotted rings, and upper layer modeling. We demonstrate in this chapter that the chosen topic is useful and interesting, and that there is much work yet to be done. We also show that there are already certain existing results which can be built upon in this dissertation.

Chapter III contains our derivation of the exhaustive limited (EL) service MAC layer model for token rings. We show that this model very accurately models a feature of token ring service not heretofore modeled. We also take into account the dependence of each token vacation time on the previous token cycle time by applying the effect of packet transmissions on a previously equilibrium state at all other stations. The model is initially developed as an M/G/1 model, and then modified so that the actual arrival process assumed in Chapter V (a non-Markov process which we will term "series" arrivals and abbreviate "S") is used, thus converting it to an S/G/1 model.

Chapter IV is the development of a priority model for exhaustive limited service. We derive this model assuming both Poisson arrivals and series arrivals. We also apply this priority scheme to existing slotted ring models. We will show that the exhaustive limited model derived in the previous chapter is identical to the priority model when the number of priorities is assumed to be equal to one.

Chapter V is our derivation of a model for the protocol stack. We present this model in a parameterized form, that is, the number of queues modeled is a parameter in the model. Our derivation assumes that the initial arrivals to the network follow a Poisson process; using this assumption, we then derive the arrival processes for each of the queues in the series, giving both their probability density functions and Laplace transforms. We then modify the model so that priority service is taken into account, which will allow us to successfully merge the upper layer stack model with the exhaustive limited service model for the MAC layer. We then show how the upper layer stack model applies to slotted rings, using existing slotted ring queueing models.

In Chapter VI we cover the simulation system used. We state the assumptions made in its development, its features, and its limitations. We also show the mathematical methods used to ensure statistical validity.

Chapter VII compares the results given by the queueing models to the results from the simulations. We initially establish a "base case" for single priority operation, as well as one for multiple priority operation, and compare the queueing model to the simulation for both of these cases. Then comparisons are made by varying each parameter of the base case in turn and indicating the results. Absolute and relative agreement are indicated, as well as confidence intervals. Comparisons are made for both token rings and slotted rings.

Recognizing the fact that few actual networks will have homogeneous traffic modes, Chapter VIII uses the simulation system to model various configurations of non-homogeneous traffic on both token rings and slotted rings. Traffic will be modeled on an end-to-end basis, thus showing the effects of the upper layers on message delay and queue lengths.

Chapter IX discusses the applicability of the analytic and simulation models to real world problems. We discuss layer groupings, and how they are mapped to the analytic model. Also discussed are the uses of the simulation system to gain prediction data that can be used by network designers.

We present our conclusions in Chapter X. We will make statements about the applicability of the models derived, as well as the simulation system when used alone. We also state what future work is to be done.

4 Preliminary Results

Preliminary results so far include a mean value analysis of EL service, and progress on the derivation of the z-transform for EL service at the MAC layer. These models are obviously important as an aid to the understanding of token ring operation. Both versions of the EL model assume a Poisson arrival process; this assumption will have to be changed to take series arrivals into account.

The mean value EL model is much easier to use, as it is a much simpler model. Its drawbacks, however, are basically twofold. First, it does not take into account the dependence of the token cycle time on the previous token cycle; second, it assumes that the variance of the token vacation time is assumed to be unchanged from its exhaustive service value. We show in Appendices A and B, however, that neither of these assumptions cause results to vary significantly from simulation results where these changes are taken into account.

In the derivation of the z-transform for the EL model, we correct the deficiencies of the approximate model, in that token vacation time dependencies are accounted for, and the variance of the token vacation time is explicitly given. It has two major drawbacks and one minor one. The first major drawback is that it is not so much a parameterized model, but a parameterized solution method; that is, there does not exist a general form for the z-transform of the distribution of queue lengths, but a parameterized algorithm for deriving the transform does exist. The second major drawback is that even though a solution theoretically exists for any limit $n \geq 1$, computational difficulties render these solutions difficult to obtain for large n (fortunately, EL service disciplines with large values of n tend to behave like exhaustive service disciplines). The minor drawback is that it is also more difficult to program than the approximate model, as each value of n has a different transform. It does, however, greatly increase the understanding of token ring EL service.

The model for the upper layers of the OSI stack is still in its formulative stage. However, certain results have been derived. The model used will be queues in series, one queue for each layer of the OSI stack being modeled. It is easily shown that the mean arrival rate for each of the queues is identical to the arrival rate at the initial queue; however, the departure distribution from each of the queues is not Markovian.

Existing results include a derivation of the LST and probability density function for the packet departure rate for an M/G/1 queue. From this, the probability density function of the departure rate from each successive queue is found. Since the density function for the arrival rate at each queue is known, it is possible to derive exact results for the system, given that the initial arrival rate to the system is Markovian.

4.1 Analysis of the EL Model

This analysis is for EL service, where the number of packets per token may be any finite number greater than or equal to one; the analysis implicitly includes a derivation for single-packet-per-token service. In this model, we are removing the simplifying assumptions made in the development of the approximate models. We now assume changes in the token vacation time and its variance, by accounting for the time dependence of each token cycle on the previous one.

4.1.1 Fundamental Relationships

The fundamental equation of our system is

$$q_{n+1} = q_n + \alpha_n - \beta_n \quad (1)$$

where

q_n = the number of packets present at the n^{th} token departure

α_n = the number of packets which arrive between the n^{th} and $n+1^{st}$ token departures

β_n = the number of packets transmitted between the n^{th} and $n+1^{st}$ token departures

This is obviously a time-dependent process. In order to remove the time dependency, we take the limit as $n \rightarrow \infty$, resulting in

$$q = q + \alpha - \beta \quad (2)$$

From this equation, we see that queue length at token departure is the natural point at which to embed a semi-Markov process. This will allow the derivation of the z-transform of the system.

The z-transform is dependent upon a finite subset of the state probabilities of the system. This subset consists of states $0 \cdots \omega - 1$, where each state represents the number of packets enqueued at token departure. The symbol ω represents the maximum number of packets per token which may be transmitted by any station. The semi-Markov chain showing the state transition probabilities for state $\epsilon \geq \omega$ is shown in Figure 3:

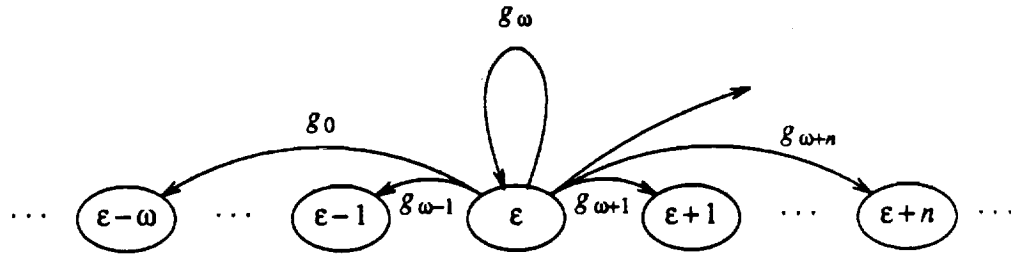


Figure 3

The transition probabilities given that the Markov chain is in state 0 are shown in Figure 4; the state probabilities given that the chain is in state 1 are shown in Figure 5.

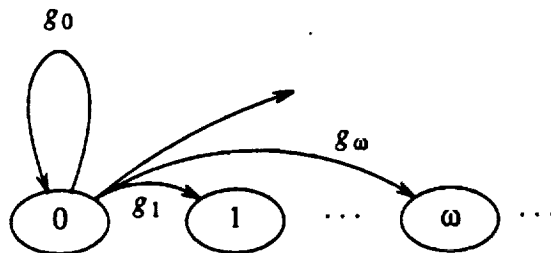


Figure 4

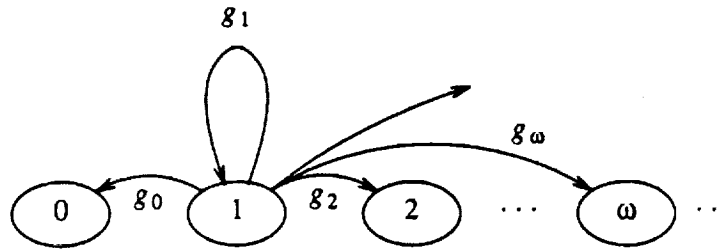


Figure 5

The probability transition matrix P representing the transition probabilities for the entire system is shown in Figure 6.

$$P = \begin{bmatrix} P_{0|0} & P_{1|0} & P_{2|0} & \cdots & P_{\omega|0} & P_{\omega+1|0} & P_{\omega+2|0} & \cdots \\ P_{0|1} & P_{1|1} & P_{2|1} & \cdots & P_{\omega|1} & P_{\omega+1|1} & P_{\omega+2|1} & \cdots \\ P_{0|2} & P_{1|2} & P_{2|2} & \cdots & P_{\omega|2} & P_{\omega+1|2} & P_{\omega+2|2} & \cdots \\ \vdots & \vdots & \vdots & \cdots & \vdots & \vdots & \vdots & \cdots \\ P_{0|\omega} & P_{1|\omega} & P_{2|\omega} & \cdots & P_{\omega|\omega} & P_{\omega+1|\omega} & P_{\omega+2|\omega} & \cdots \\ 0 & P_{0|\omega} & P_{1|\omega} & \cdots & P_{\omega-1|\omega} & P_{\omega|\omega} & P_{\omega+1|\omega} & \cdots \\ 0 & 0 & P_{0|\omega} & \cdots & P_{\omega-2|\omega} & P_{\omega-1|\omega} & P_{\omega|\omega} & \cdots \\ \vdots & \vdots & \vdots & \cdots & \vdots & \vdots & \vdots & \cdots \end{bmatrix}$$

Figure 6

Each entry P_{ij} is defined as

$$P_{ij} \equiv P_{j|i} = P[\epsilon_{t=n+1} = j \mid \epsilon_{t=n} = i] \quad (3)$$

The computation of each $P_{j|i}$ requires the evaluation of packet arrival probabilities. For any state $\epsilon < \omega$, these probabilities break down into two distinct cases. The first case is transitions to state $\epsilon = 0$; the second is transitions to all other states. Defining the probability of k arrivals given state ϵ to be $\pi_{k|\epsilon}$, the probability of a transition to state $\epsilon = 0$ is given by

$$P_{0|i} = \sum_{k=0}^{\omega-i} \pi_{k|i} \quad (4)$$

The probability of a transition to state $\epsilon > 0$ is given by

$$P_{j|i} = \pi_{\omega-i+j|i} \quad (5)$$

The calculation of each $\pi_{j|i}$ is complicated by the fact that the arrivals are constrained to arrive early enough to make a contribution to service time. For example, it is evident that for any packets to be left behind at token departure, the maximum number of packets per token allowed must have been transmitted, which we assume for this example is equal to three (if less than ω packets were delivered, the station was empty of packets at token departure). However, if only a single packet was present at the previous token departure, and the number left behind at the current token departure is greater, then at least four packets must have arrived during a total of three packet service times. But if all three arrived during the third service time, there would only have been a single service time, (a contradiction) since only the single packet present would have been served with the token then departing. Therefore, some of the packets had to have arrived prior to the third service time for there to have been ω packets delivered. An example of this process can be seen in Figure 7.

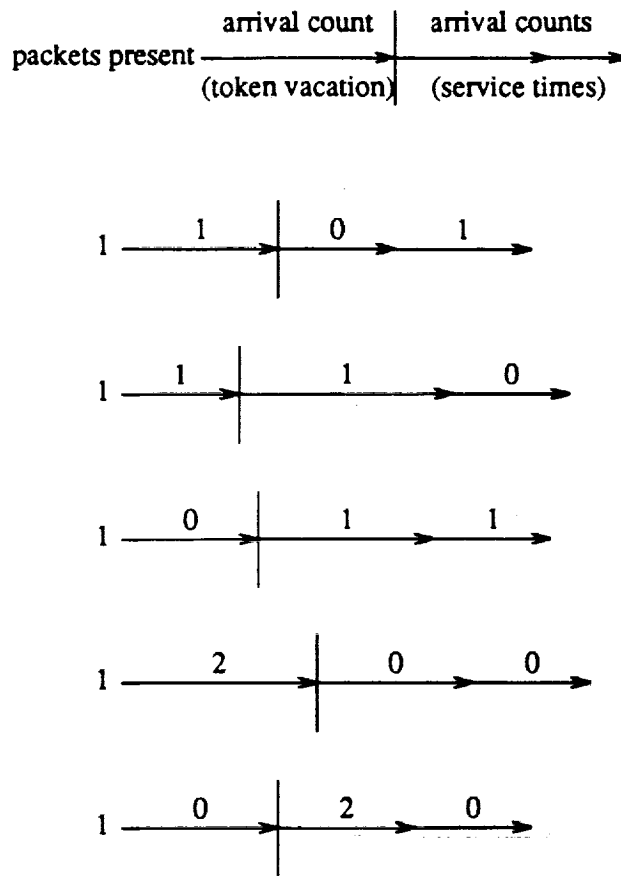


Figure 7

To account for the above constraints on arrivals, the computation of the probability of k arrivals is done by "splitting" the possible arrival time frame into units consisting of the token vacation time and the appropriate number of packet service times. Arrival probabilities are then calculated for the proper number of arrivals split among the various time frames, and multiplied by an inclusion function. This inclusion function takes on the values 0 or 1 depending on whether or not an arrival is "early enough." The equation for $\pi_{i,j}$ is

$$\pi_{i,j} = \sum_{k_0=0}^i P[\alpha_0=k_0] I(0, \sum_{x=0}^0 k_x, j) \sum_{k_1=0}^{i-\sum_{x=0}^0 k_x} P[\alpha_1=k_1] I(1, \sum_{x=0}^1 k_x, j) \sum_{k_2=0}^{i-\sum_{x=0}^1 k_x} P[\alpha_2=k_2] I(2, \sum_{x=0}^2 k_x, j)$$

$$\begin{aligned}
& \dots \sum_{k_{\omega}=0}^{i-\sum_{x=0}^{\omega-1} k_x} P[\alpha_{\omega}=k_{\omega}] I(\omega, \sum_{x=0}^{\omega} k_x, j) \\
& = \prod_{s=0}^{\omega} \sum_{k_s=0}^{i-\sum_{x=0}^{s-1} k_x} P[\alpha_s=k_s] I(s, \sum_{x=0}^s k_x, j)
\end{aligned} \tag{6}$$

The inclusion function $I(a, b, j)$ is defined by Equation (7). The quantity a is the "slot number," with the token vacation time being counted as slot 0, the initial packet service time as slot 1, etc., the quantity b is the number of total arrivals so far, and the quantity j is the number of packets left behind at the previous token departure.

$$\begin{aligned}
I(a, b, j) &= 1, \quad j + b \geq a \\
I(a, b) &= 0, \quad \text{otherwise}
\end{aligned} \tag{7}$$

4.1.2 The z-transform

The method we use to derive the z-transform is to derive the conditional transforms which will thus allow us to express the unconditional transforms in terms of state probabilities. For any ω , there will be exactly ω unknown probabilities. Defining $B^*(s)$ as the LST of the service time distribution for a single packet, and $V^*(\lambda - \lambda z | S)$ as the LST of the token vacation time given state S , the z-transform given state $\epsilon \geq \omega$ is

$$G(z | \epsilon \geq \omega) = V^*(\lambda - \lambda z | \epsilon \geq \omega) [B^*(\lambda - \lambda z)]^{\omega} z^{\epsilon - \omega} \tag{8}$$

where

$$\begin{aligned}
G(z) &\equiv \sum_{k=0}^{\infty} z^k P[X=k] \\
B^*(s) &\equiv \int_{-\infty}^{\infty} e^{-st} b(x) dx
\end{aligned}$$

For any state $0 \leq \epsilon < \omega$, the calculation of the z-transform is more complex. It is a multistage process whereby the final z-transform is constructed from intermediate transforms. The intermediate transforms are added to each other after multiplying each intermediate transform by the appropriate probability, thus ensuring that our final transform is of a valid probability function. The calculation of the z-transform for state $\epsilon = 1$, with $\omega = 2$ follows.

The initial step is to write down the transition probabilities in terms of arrival probabilities:

$$P_{0|1} = \pi_{0|1} + \pi_{1|1}$$

$$P_{1|1} = \pi_{2|1}$$

$$P_{2|1} = \pi_{3|1} \tag{9}$$

etc., where

$P_{x|y}$ \equiv probability of a transition to state x given state y

$\pi_{x|y}$ \equiv probability of x arrivals given state y

From the above series we obtain two transforms which we denote $G_0(z|1)$, and $G_1(z|1)$. They are obtained by solving

$$G_0(z|1) = \frac{1}{1 - \pi_{1|1}} \left[\pi_{0|1} z^0 + \pi_{2|1} z^1 + \pi_{3|1} z^2 + \dots \right] \tag{10}$$

and

$$G_1(z|1) = \frac{1}{1 - \pi_{0|1}} \left[\pi_{1|1} z^0 + \pi_{2|1} z^1 + \pi_{3|1} z^2 + \dots \right] \tag{11}$$

respectively, which result in

$$G_0(z|1) = \frac{\pi_{011} + z(\pi_{111} - \pi_{011})}{1 - z(1 - \pi_{111})} \quad (12)$$

and

$$G_1(z|1) = \frac{\pi_{010}}{1 - z(1 - \pi_{011})} \quad (13)$$

Defining $G(z|1)$ to be the conditional z -transform given state $\epsilon = 1$, $G(z|1)$ is obtained by adding $k_0 G_0(z|1)$ and $k_1 G_1(z|1)$. It is immediately obvious that

$$k_0 = \frac{1 - \pi_{111}}{2 - \pi_{011} - \pi_{111}} \quad (14)$$

and

$$k_1 = \frac{1 - \pi_{011}}{2 - \pi_{011} - \pi_{111}} \quad (15)$$

Each of the conditional transforms are derived in a similar manner.

Once all the conditional transforms have been obtained, the final unconditional transform $G(z)$ for $\omega = 2$ is given by

$$G(z) = P_0 G(z|0) + P_1 G(z|1) + \sum_{\epsilon=2}^{\infty} P_{\epsilon} G(z|\epsilon) \quad (16)$$

This ultimately resolves to

$$G(z) = \frac{V^*(\lambda - \lambda z | \epsilon \geq 2) [B^*(\lambda - \lambda z)]^2 \sum_{\epsilon=0}^1 P_{\epsilon} z^{\epsilon} - z^2 \sum_{\epsilon=0}^1 P_{\epsilon} G(z|\epsilon)}{V^*(\lambda - \lambda z | \epsilon \geq 2) [B^*(\lambda - \lambda z)]^2 - z^2} \quad (17)$$

It is evident that the z -transform given any $\omega < \infty$ takes the following general form:

$$G(z) = \frac{V^*(\lambda - \lambda z | \varepsilon \geq \omega) [B^*(\lambda - \lambda z)]^\omega \sum_{\varepsilon=0}^{\omega-1} P_\varepsilon z^\varepsilon - z^\omega \sum_{\varepsilon=0}^{\omega-1} P_\varepsilon G(z | \varepsilon)}{V^*(\lambda - \lambda z | \varepsilon \geq \omega) [B^*(\lambda - \lambda z)]^\omega - z^\omega} \quad (18)$$

Note however that the form for each $G(z | \varepsilon)$ is dependent on ω , thus requiring a separate solution for each ω limit.

The unknown state probabilities $P_0, P_1, \dots, P_\omega$ are obtained using the method in [GENT88]. The unconditional transform has a degree ω polynomial in z in both the numerator and the denominator. Since the transform is analytic, Rouché's Theorem guarantees that the numerator and denominator will have the same roots. The denominator polynomial, which involves only z , is used to obtain the zeros of the numerator polynomial, resulting in ω equations in ω unknowns. The necessary $(\omega+1)^{th}$ equation is obtained using the fact that $G(z) \big|_{z=1} = 1$. These zeros are now used in the numerator to solve for $P[0], P[1], \dots, P[\omega-1], P[\omega]$.

4.1.3 Token Cycle Time

Using $E[W]$ to represent the expected waiting time of a packet, $E[Q]$ to represent the expected queueing delay, and Γ_v to be the random variable representing the token vacation time distribution, it is shown in [LEVY75] that $E[W]$ is given by:

$$E[W] = E[Q] + \frac{E[\Gamma_v^2]}{2E[\Gamma_v]} \quad (19)$$

The quantity $E[Q]$ is given by the well-known Pollaczek-Kinchine mean-value equation

$$E[Q] = \rho + \rho^2 \frac{1 + C^2}{2(1 - \rho)} \quad (20)$$

where C^2 is the square of the coefficient of variation of the service time, and ρ is the so-called utilization factor at a single station.

The mean token cycle time $E[\Gamma]$ (Γ is the token cycle time random variable) is unaffected by the service mode so long as the network operation is stable. $E[\Gamma]$ is given by

$$E[\Gamma] = \frac{\gamma}{1 - N\rho} \quad (21)$$

where N is the number of ring stations, and γ is the ring latency, also known as the empty ring walk time. The unknown quantity in Equation (15) is $E[\Gamma_v^2]$, the second moment of the token vacation time random variable.

4.1.4 Token Vacation Time Second Moment

To obtain $E[\Gamma_v^2]$, we make use of the derivation of the variance of the token vacation time in [KONH74]. Konheim and Meister give this variance as

$$\text{Var}[\Gamma_v] = \frac{\gamma(N-1)\text{Var}[L]}{(1 - N\rho)^2} \quad (22)$$

where L is the station load random variable. $\text{Var}[L]$ is given by

$$\text{Var}[L] = E[S^2]E[A^2] - \rho^2 \quad (23)$$

where S is the service time random variable, and A is the arrival rate random variable. The variance of Γ_v can also be expressed as

$$\text{Var}[\Gamma_v] = \frac{(N-1)E^2[\Gamma_v]E[L^2]\mu}{\rho\lambda\gamma} \quad (24)$$

which proves that the variance of Γ_v is inversely proportional to the product of the traffic intensity at a single station and the expected number of arrivals in an empty ring walk time (the mean number of arrivals in the shortest possible token cycle).

We now present a corollary to the proof in [KONH74]. The mean number of packets per token actually transmitted is unaffected by the maximum number of packets per token allowed.

thus rendering Konheim and Meister's proof valid for any limited service discipline. Therefore, since it is shown above that the variance of the token vacation time is inversely proportional to the mean number of packet arrivals in a minimum token vacation time, the limited service token vacation time variance for a limited service discipline is the product of the traffic intensity at a single station and the minimum number of packets present at the station after one token vacation time. This is given by

$$Var[\Gamma_v] = \frac{(N-1)E^2[\Gamma_v]E[L^2]\mu}{\rho(\lambda\gamma + E[R])} \quad (25)$$

Using the fact that $E[X^2] = E^2[X] + Var[X]$, the expression for the second moment of the token vacation time $E[\Gamma_v^2]$ is

$$E[\Gamma_v^2] = \frac{(N-1)E^2[B]}{\lambda\gamma + E[R]} + \frac{\gamma^2(1-\rho)^2}{(1-N\rho)^2} \quad (26)$$

In order to find $E[R]$ it is only necessary to note that the state of the system at token departure is exactly this quantity. The second moment of the token cycle time can now be found by a simple substitution.

4.2 Results for Queues in Series

Given that the arrivals at a queue are Poisson, the LST $D^*(s)$ of the probability density function and the density function $D(t)$ of the time between departures are

$$A^*(s) = \frac{\lambda - \lambda(1-\rho)\Theta^*(s)}{s} \quad (27)$$

$$a(t) = \lambda - \lambda(1-\rho)\Theta(t) \quad (28)$$

where

$f(t)$ = the probability density function of the interarrival process

$b(t)$ = the probability density function of the service time distribution

$\theta(t) \equiv$ the convolution of $f(s)$ and $b(s)$

It is evident that the density function for the departures from a queue is the density function for arrivals at the next queue in the series, so we immediately have the interarrival time densities.

These results are used to find the LST and probability density function for the interarrival time densities to each of the successive queues. If we designate the initial exponential interarrival time density as $a_1(t)$, then the LST of the interarrival time density $A_k^*(s)$ and the density itself $a_k(t)$ for each queue, where $k = 2, 3, \dots, n$ are

$$A_k^*(s) = \frac{\lambda - \lambda(1-\rho)\Theta_{k-1}^*(s)}{s} \quad (29)$$

$$a_k(t) = \lambda - (1-\rho)\Theta_{k-1}(t) \quad (30)$$

The above gives the density and LST for the packet interarrival time at each queue. We also need to find the probability density function for the number of arrivals in a set period of time. This density will be a conditional density, as we are not dealing with a memoryless distribution.

The density function for the probability of k arrivals in time t is derived as follows. It is evident that since the density function of the time until the next arrival is not memoryless, the probability of an arrival during time t will depend on the number of arrivals that have already occurred, and also on when they occurred. This results in the k -fold convolution of the density function with itself (defining the one-fold convolution to be the density function, and the zero-fold convolution to be the probability that zero arrivals occur during time t). Therefore, given time t_0 , the probability of k i^{th} level series arrivals during time $t - t_0$ is

$$P[k | t_0] = a_i(t - t_0)^{[k]} \quad (31)$$

where $x^{[y]}$ denotes the y -fold convolution of the density function $x(\cdot)$ with itself.

The above are the results derived to date for series queues. However, it is evident from the work done so far that useful results can be derived for bulk arrivals and priority traffic. It will then be necessary to combine the model for queues in series with the MAC layer models for both token rings and slotted rings. The method employed here will be to modify the arrival processes assumed in the MAC layer models so that they conform to the departure processes from the model for queues in series. From there, it is a relatively simple matter to find mean delay, queue lengths, and the total number of packets in the system (as well as the total number of messages, which may be quite different).

5 Bibliography

- [APPL86] J. M. Appleton and M. M. Peterson, Traffic Analysis of a Token Ring PBX, *IEEE Transactions on Communications COM-34*, 5 (May 1986), 417-422.
- [BHUY89] L. N. Buuyan, D Ghosal and Q. Yang, Approximate Analysis of Single and Multiple Ring Networks, *IEEE Transactions on Computers* 38, 7 (July 1989), 1027-1040.
- [EVER89] D. Everitt, A Note on the Pseudoconservation Laws for Cyclic Service Systems with Limited Service Disciplines, *IEEE Transactions on Communications COM-37*, 7 (July 1989), 781-783.
- [GENT88] W. L. Genter and K. S. Vastola, Performance Measurement and Analysis of the Token Bus Network, *Proceedings of the 27th IEEE Conference on Decision and Control*, Austin, Texas, Dec. 7-9, 1988, 1489-1495.
- [GORU87] M. R. Gorur, *Priorities and Dynamic Ring Membership in an 802.4 Network*, Master's Thesis, University of Virginia, Charlottesville, Virginia, 1987.
- [HEYM83] D. P. Heyman, Data-Transport Analysis of Fasnnet, *Bell System Tech. J.* 62, 8 (Oct. 1983), 2547-2560.
- [KONH72] A. G. Konheim and B. W. Meister, Waiting Lines and Times in a System with Polling, *J. ACM* 21, (July 1974), 470-490.
- [KUEH79] P. J. Kuehn, Multiqueue Systems with Nonexhaustive Cyclic Service, *Bell System Tech. J.* 58, 3 (Mar. 1979), 671-698.
- [LEVY75] Y. Levy and U. Yechiali, Utilization of Idle Time in an M/G/1 Queueing System, *Management Science* 22, 2 (Oct. 1975), 202-211.
- [LI88] S. Li and M. El Zarka, Dynamic Bandwidth Allocation on a Slotted Ring with Integrated Services, *IEEE Transactions on Communications COM-36*, 7 (July 1988), 826-833.
- [MANF85] D. R. Manfield, Analysis of a Priority Polling System for Two-Way Traffic, *IEEE Transactions on Communications COM-33*, 9 (Sep. 1985), 1001-1006.
- [MILL87] K. Mills, M. Wheatley and S. Heatley, Prediction of Transport Layer Performance Through Simulation, *Proceedings of the Workshop on Factory Communications*, National Bureau of Standards, Gaithersburg, Maryland, March 17-18, 1987.
- [MITC86] L. C. Mitchell and D. A. Lide, End-to-End Performance Modeling of Local Area Networks, *IEEE Journal on Selected Areas in Communications SAC-4*, 6 (Sep. 1986), 975-985.
- [MURA88] M. Murata and H. Takagi, Two-Layer Modeling for Local Area Networks, *IEEE Transactions on Communications Com-36*, 9 (Sep. 1988), 1022-1034.
- [PEDE87] J. Peden and A. Weaver, Performance of Priorities on an 802.5 Token Ring,

Proceedings of Sigcomm 1987 Workshop: Frontiers in Computer Communications Technology, Stowe, Vermont, Aug. 11-13, 1987, 58-66.

- [PEDE88] J. H. Peden and A. C. Weaver, The Utilization of Priorities on Token Ring Networks, *Proceedings of the 13th Conference on Local Computer Networks*, Minneapolis, Minn., Oct. 10-12, 1988, 472-478.
- [REGO84] V. J. Rego and L. M. Ni, *Performance Modeling of Token-Passing Protocols*, Michigan State University Technical Report, 1984.
- [SCHU89] N. Schult, *Analytic Models for Token-Ring Networks*, Ph.D. Dissertation, University of Virginia, Charlottesville, Virginia, May, 1989.
- [STRA88] W. T. Strayer, *Performance Analysis of the Manufacturing Automation Protocol*, Master's Thesis, University of Virginia, Charlottesville, Virginia, 1988.
- [TAKA86] H. Takagi and M. Murata, Queueing Analysis of Nonpreemptive Reservation Priority Discipline, *Proceedings of Performance '86 and ACM Sigmetrics (Joint Conference)*, Raleigh, NC, May 27-30, 1986, 237-241.
- [ZAFI87] M. Zafirovic-Vukotic and I. G. Neimegeers, Performance Modelling of the Orwell Basic Access Mechanism, *Proceedings of the SigComm '87 Workshop on Frontiers in Computer Communications Technology*, Stowe, Vermont, August 11-13, 1987, 35-48.